

Neural Text deGeneration

Inconsistency & Unlikelihood

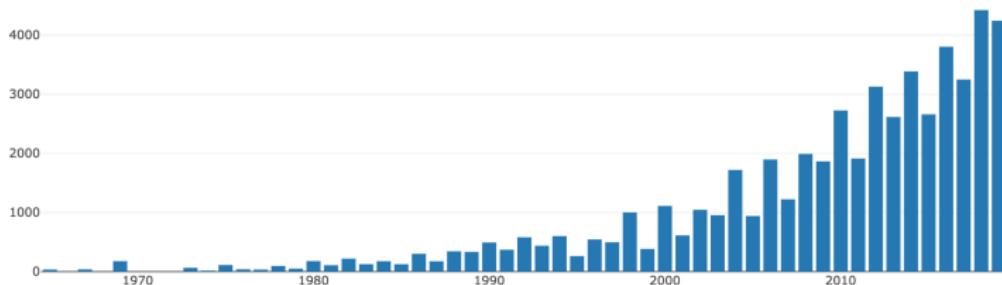
Sean Welleck

NYU

18th March, 2020

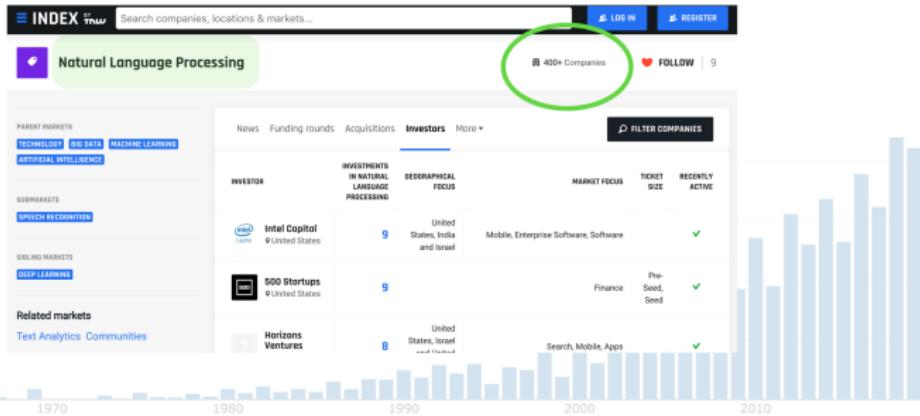
- ▶ Despite substantial progress...

Yearly Paper Distribution



nlpexplorer.org

► Despite substantial progress...



- ▶ Despite substantial progress...

Better Language Models and Their Implications

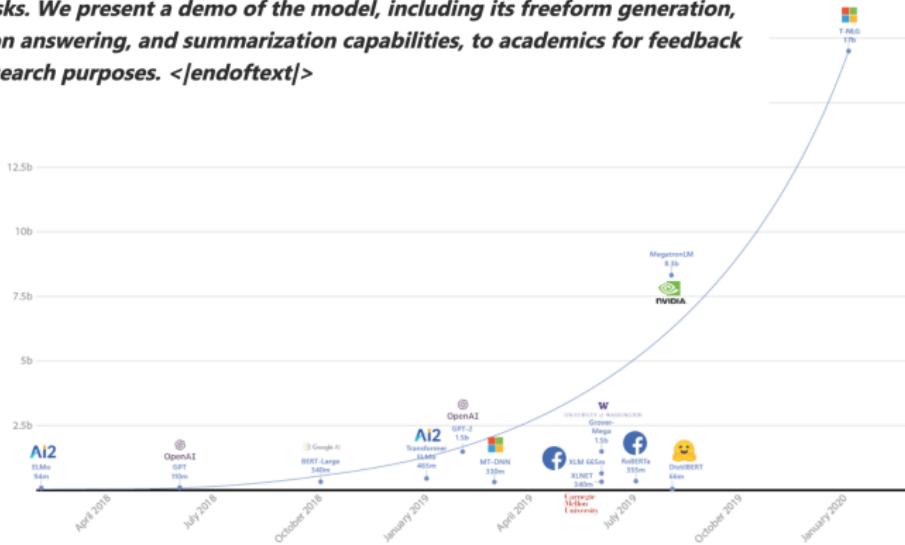
We've trained a large-scale unsupervised language model which generates coherent paragraphs of text, achieves state-of-the-art performance on many language modeling benchmarks, and performs rudimentary reading comprehension, machine translation, question answering, and summarization—all without task-specific training.

Due to concerns about large language models being used to generate deceptive, biased, or abusive language at scale, we are only releasing a much smaller version of GPT-2 along with sampling code. We are not releasing the dataset, training code, or GPT-2 model weights. Nearly a year ago we wrote in the OpenAI Charter: “we

<https://openai.com/blog/better-language-models/>

- ▶ Despite substantial progress...

Turing Natural Language Generation (T-NLG) is a 17 billion parameter language model by Microsoft that outperforms the state of the art on many downstream NLP tasks. We present a demo of the model, including its freeform generation, question answering, and summarization capabilities, to academics for feedback and research purposes. </endoftext>

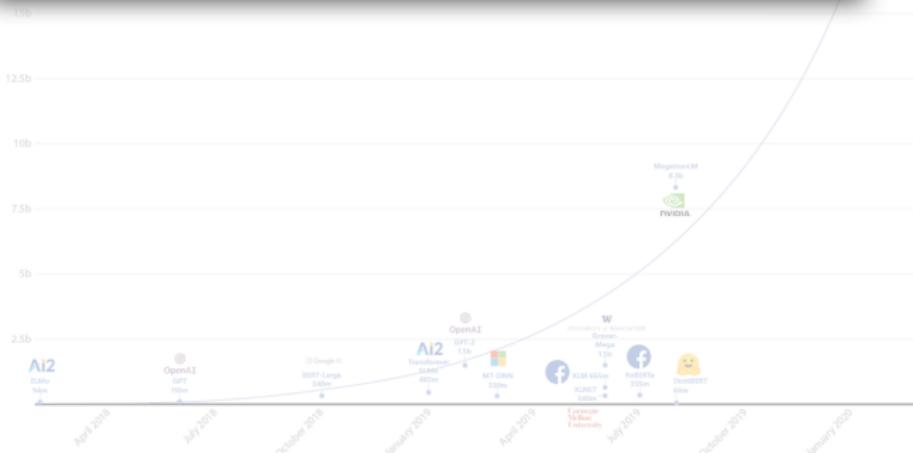


► Despite substantial progress...

Turing Natural Language Generation (T-NLG) is a 17 billion parameter language model by Microsoft that outperforms the state of the art on many downstream NLP tasks. We present a demo of the model, including its freeform generation

questions and responses, and its ability to generate text from images.

- This summary was generated by the Turing-NLG language model itself.



- ▶ Standard neural language models show degeneration

Context:

In a shocking finding, scientist discovered a herd of unicorns living in a remote, previously unexplored valley, in the Andes Mountains. Even more surprising to the researchers was the fact that the unicorns spoke perfect English.

Continuation (BeamSearch, b=10):

GPT-2 Continuations [from Holtzmann et al ICLR 2020]

► Standard neural language models show degeneration

Prefix
 \mathcal{L}_{MLE}

... starboard engines and was going to crash . “ We 're going in , ”
he said . “ We 're going to crash . We 're going to

Repetition

Model: 40% 4-gram repetition

Human: ~1% 4-gram repetition

Large-Scale Transformer LM [from Welleck et al ICLR 2020]

This talk

- ▶ What do we mean by ‘neural text degeneration’?
(Background)
- ▶ What do we theoretically know about it?
(Inconsistency)
- ▶ What are some ways of preventing it?
(Inconsistency, Unlikelihood)

Background – Neural Sequence Modeling

Goal:

1. estimate $p_\theta(\underbrace{y_1, \dots, y_T}_{\text{sequence}} | \underbrace{C}_{\text{context}})$

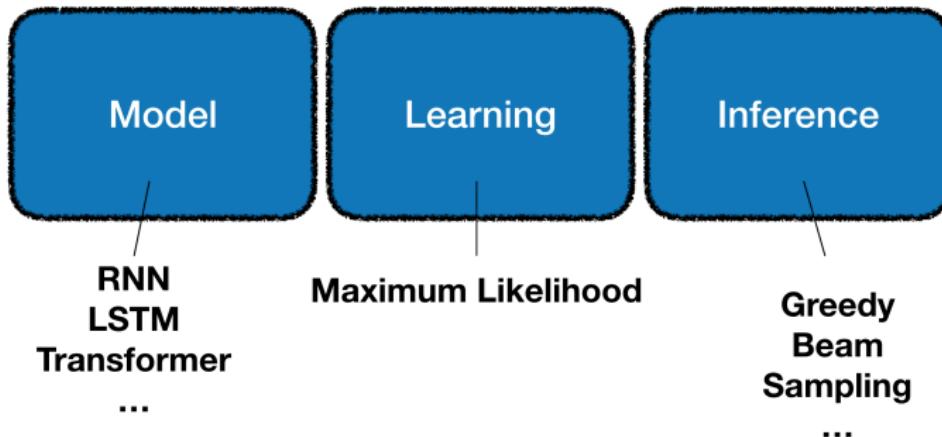
► $p_\theta(\text{how, are, you, ?} | \text{元気ですか?})$

2. decode $(\hat{y}_1, \dots, \hat{y}_T) \sim \underbrace{\mathcal{F}(p_\theta, C)}_{\text{decoding algorithm}}$

► $(\text{how, are, you, ?}) \sim \mathcal{F}(p_\theta, \text{元気ですか?})$

Machine translation, dialogue modeling, language modeling, story generation, text completion, etc.

Background – Neural Sequence Modeling

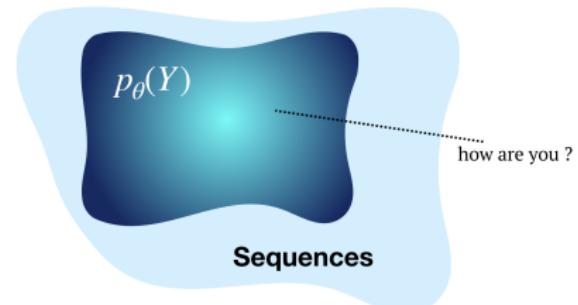


Model

- ▶ Neural autoregressive sequence model:

$$p_{\theta}(y_1, \dots, y_T | C) = \prod_{t=1}^T \underbrace{p_{\theta}(y_t | y_{<t}, C)}_{\text{neural network}}$$

- ▶ Neural network:
 - ▶ RNN / LSTM
 - ▶ Transformer



Learning

Maximum Likelihood Estimation (MLE)

- Given dataset $\{(C^{(n)}, Y^{(n)})\}_{n=1}^N$,

$$\arg \max_{\theta} \sum_{n=1}^N \sum_{t=1}^{T^{(n)}} \log p_{\theta}(y_t^{(n)} | y_{<t}^{(n)}, C^{(n)})$$

Inference

- Given trained model p_θ , context C

$$\hat{Y} = \arg \max_{Y \in \mathcal{Y}} \log p_\theta(Y|C)$$

Inference

- Given trained model p_θ , context C

$$\hat{Y} = \arg \max_{Y \in \mathcal{Y}} \log p_\theta(Y|C)$$

- Approximate using a decoding algorithm $\mathcal{F}(p_\theta, C)$
 - ▶ Greedy search
 - ▶ Beam search
 - ▶ Sampling

Inference

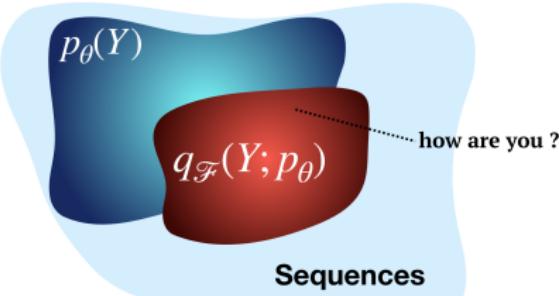
- Given trained model p_θ , context C

$$\hat{Y} = \arg \max_{Y \in \mathcal{Y}} \log p_\theta(Y|C)$$

- Approximate using a decoding algorithm $\mathcal{F}(p_\theta, C)$

- ▶ Greedy search
- ▶ Beam search
- ▶ Sampling

- The decoding algorithm induces a new sequence distribution



Inference – Decoding Algorithms

Greedy Search

Recursively selects the most likely token until `<eos>`:

$$\tilde{y}_t = \arg \max_{v \in V} \log p_\theta(y_t = v | \tilde{y}_{<t}, C).$$

Inference – Decoding Algorithms

Greedy Search

Recursively selects the most likely token until $\langle \text{eos} \rangle$:

$$\tilde{y}_t = \arg \max_{v \in V} \log p_\theta(y_t = v | \tilde{y}_{<t}, C).$$

Beam Search

Searches with a set of k high scoring prefixes:

$$P_t = \bigcup_{\rho \in P_{t-1}^{\text{top}}} \{\rho \circ v \mid v \in V\}$$

$$P_t^{\text{top}} = \arg \max_{\rho \in P_t} s(\rho).$$

Inference – Decoding Algorithms

Ancestral Sampling

$$\tilde{y}_t \sim p_{\theta}(y_t | \tilde{y}_{<t}, C).$$



Inference – Decoding Algorithms

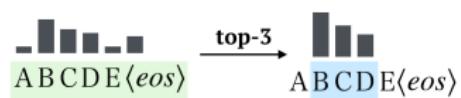
Ancestral Sampling

$$\tilde{y}_t \sim p_{\theta}(y_t | \tilde{y}_{<t}, C).$$



Top-k Sampling

$$q(v) \propto \begin{cases} p_{\theta}(v | y_{<t}, C), & v \in \arg \text{top-}k \ p_{\theta}, \\ 0, & \text{otherwise.} \end{cases}$$



Inference – Decoding Algorithms

Ancestral Sampling

$$\tilde{y}_t \sim p_{\theta}(y_t | \tilde{y}_{<t}, C).$$



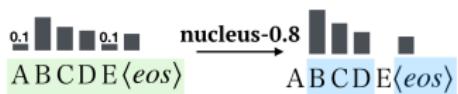
Top-k Sampling

$$q(v) \propto \begin{cases} p_{\theta}(v | y_{<t}, C), & v \in \text{arg top-k } p_{\theta}, \\ 0, & \text{otherwise.} \end{cases}$$



Top-p ('nucleus') Sampling

$$q(v) \propto \begin{cases} p_{\theta}(v | y_{<t}, C), & v \in V_{\mu}, \\ 0, & \text{otherwise.} \end{cases}$$

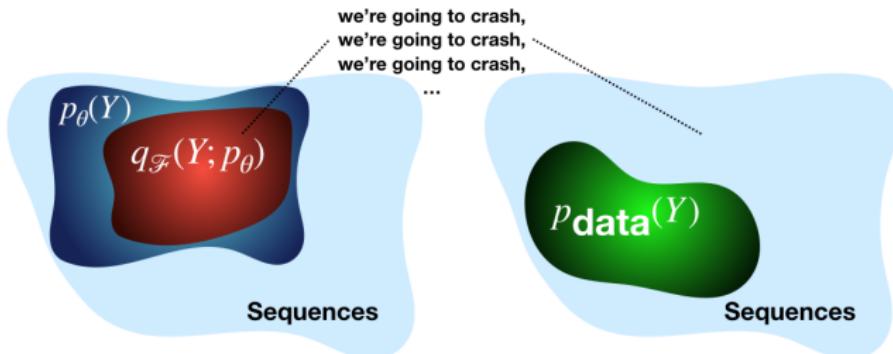


V_{μ} : smallest set with total prob $< \mu$

Problem: ‘Text Degeneration’

- ▶ Language models trained with **maximum likelihood** produce text which **differs from the true distribution** when using **common decoding algorithms**:

$$q_{\mathcal{F}}(Y) \neq p_{\text{data}}(Y)$$



Symptoms

► Repetition

Prefix
 \mathcal{L}_{MLE}

... *starboard engines and was going to crash*. "We're going in," he said. "We're going to crash. We're going to

Repetition

Model: 40% 4-gram repetition
Human: ~1% 4-gram repetition

Large-Scale Transformer LM [from Welleck et al ICLR 2020]

Symptoms

- ▶ Repetition
- ▶ Non-terminating sequences

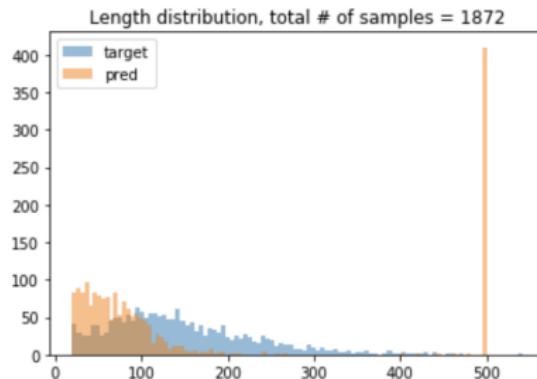


Figure: GPT-2 Completions (Wikitext-103 fine-tuned)

Symptoms

- ▶ Repetition
- ▶ Non-terminating sequences
- ▶ Logical contradictions

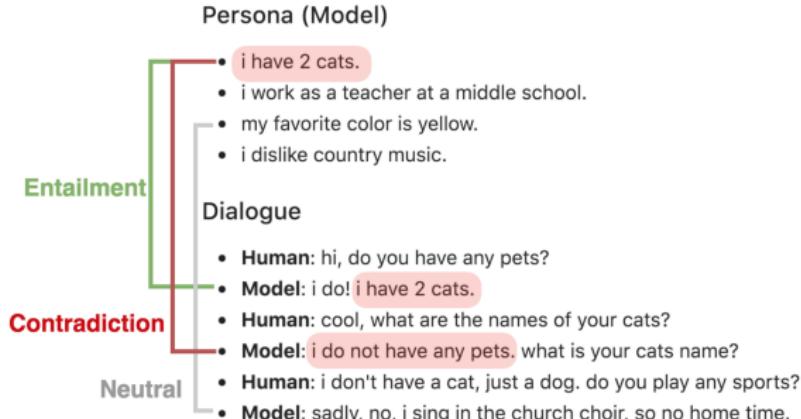


Figure: Welleck et al 2019, Dialogue NLI

Symptoms

- ▶ Repetition
- ▶ Non-terminating sequences
- ▶ Logical contradictions

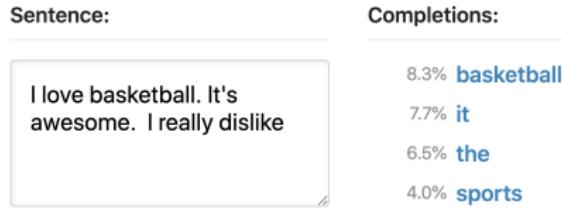


Figure: GPT-2 completion

Consistency of a Recurrent Language Model with Respect to Incomplete Decoding

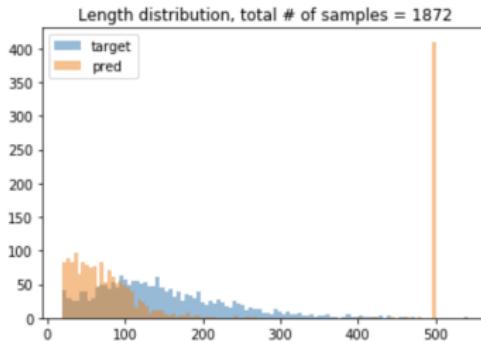
Welleck, Kulikov, Kim, Pang, Cho 2020 (under review)



Consistency of a Recurrent Language Model

- ▶ Completions appear to be *infinite-length*

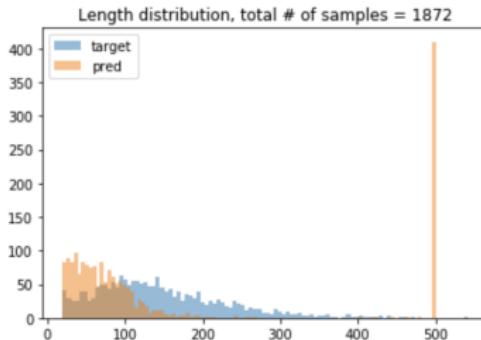
Prefix ... *starboard engines and was going to crash . “ We 're going in ,”*
Completion he said . “ We 're going to crash . We 're going to ...



Consistency of a Recurrent Language Model

- ▶ Completions appear to be *infinite-length*

Prefix ... *starboard engines and was going to crash . “ We 're going in ,*
Completion *he said . “ We 're going to crash . We 're going to ...*



- ▶ Can we *formally describe* and *analyze* what's happening?
- ▶ Can we fix it?

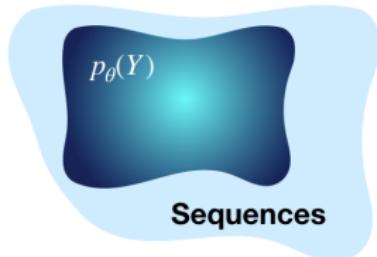
► Recurrent Language Model

Definition 2.3 (Recurrent language model). A recurrent language model p_θ is a neural network that computes the following conditional probability at each time step

$$p_\theta(y_t = v \mid y_{<t}, C) = \frac{\exp(u_v^\top h_t + c_v)}{\sum_{v' \in V} \exp(u_{v'}^\top h_t + c_{v'})},$$

$$p_\theta(Y \mid C) = \prod_{t=1}^T p_\theta(y_t \mid y_{<t}, C),$$

dog whose name is spot . i have a pet
context

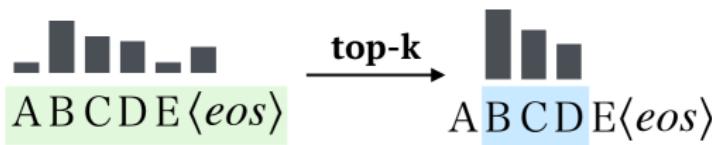


► Incomplete Decoding Algorithm

Only considers a subset of tokens at each time step

Definition 2.11 (Incomplete Decoding). *A decoding algorithm \mathcal{F} is incomplete when for each context C and prefix $y_{<t}$, there is a strict subset $V'_t \subsetneq V$ such that*

$$\sum_{v \in V'_t} q_{\mathcal{F}}(y_t = v \mid y_{<t}, C) = 1.$$

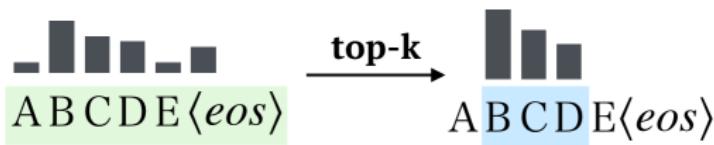


► Incomplete Decoding Algorithm

Only considers a subset of tokens at each time step

Definition 2.11 (Incomplete Decoding). *A decoding algorithm \mathcal{F} is incomplete when for each context C and prefix $y_{<t}$, there is a strict subset $V'_t \subsetneq V$ such that*

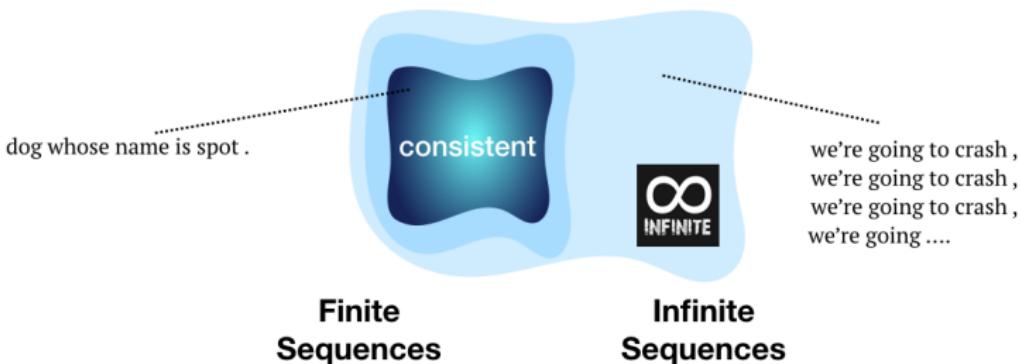
$$\sum_{v \in V'_t} q_{\mathcal{F}}(y_t = v \mid y_{<t}, C) = 1.$$



- Greedy, beam, top-k, nucleus,...

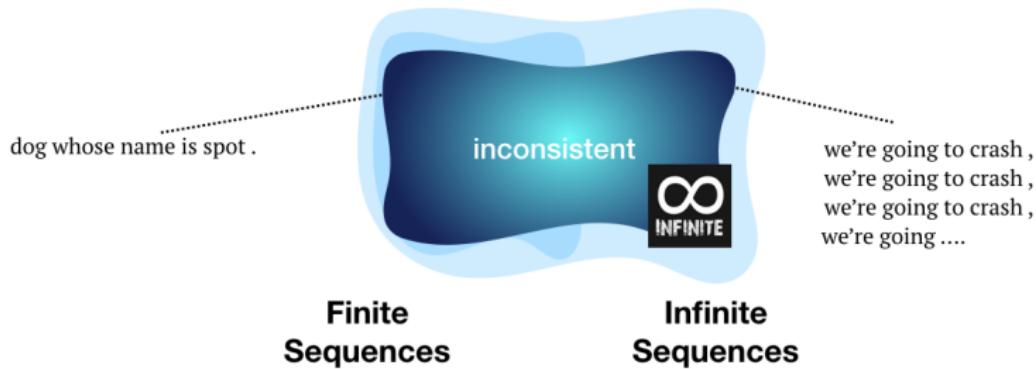
► Consistency

Definition 3.1 (Consistency of a recurrent language model).
A recurrent language model is consistent under a context distribution $p(C)$ if $p_\theta(|Y| = \infty) = 0$. Otherwise, the recurrent language model is said to be inconsistent.



► Consistency

Definition 3.1 (Consistency of a recurrent language model).
A recurrent language model is consistent under a context distribution $p(C)$ if $p_\theta(|Y| = \infty) = 0$. Otherwise, the recurrent language model is said to be inconsistent.



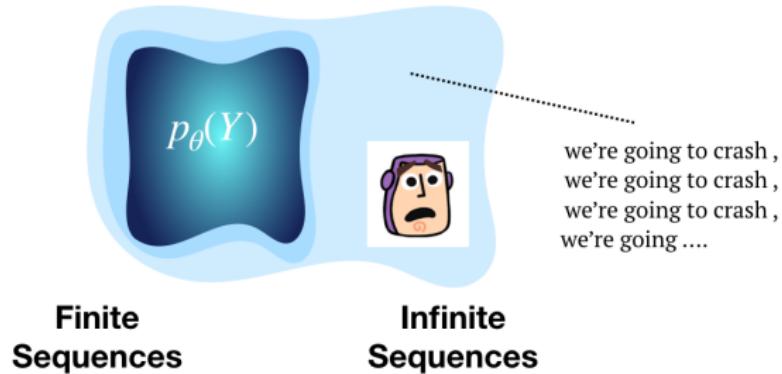
- ▶ Is the **model** consistent?¹

¹See also [Chen et al 2017, Recurrent Neural Networks as Weighted Language Recognizers]

- Is the **model** consistent?¹ Yes!

Lemma 3.2. *A recurrent language model p_θ is consistent if $\|h_t\|_p$ is uniformly bounded for some $p \geq 1$.*

Tanh, layer normalization, ...

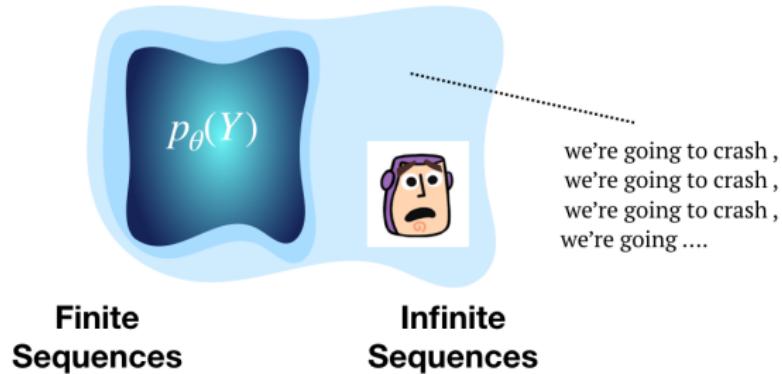


¹ See also [Chen et al 2017, Recurrent Neural Networks as Weighted Language Recognizers]

- ▶ Is the **model** consistent?¹ Yes!

Lemma 3.2. *A recurrent language model p_θ is consistent if $\|h_t\|_p$ is uniformly bounded for some $p \geq 1$.*

Tanh, layer normalization, ...



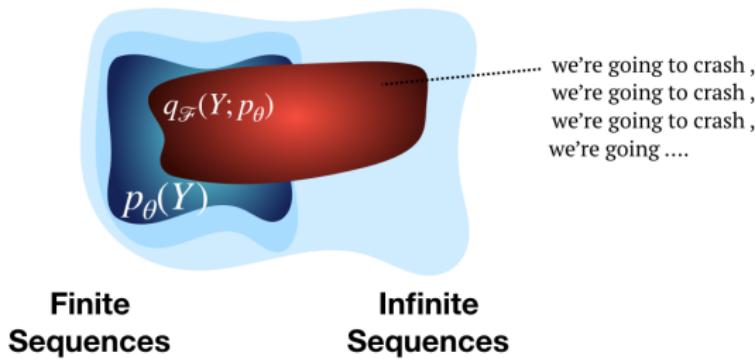
- ▶ $p_\theta(\text{infinite sequence}) = 0$

¹ See also [Chen et al 2017, Recurrent Neural Networks as Weighted Language Recognizers]

- ▶ Is the **model + decoding algorithm** consistent?

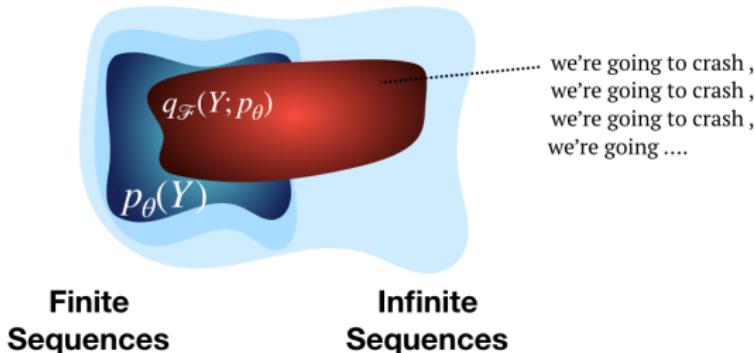
- Is the **model + decoding algorithm** consistent? **No!**

Theorem 3.4 (Inconsistency) of an incomplete decoding algorithm). *There exists a consistent recurrent language model p_θ from which an incomplete decoding algorithm \mathcal{F} , that considers only up to $(|V| - 1)$ -most likely tokens according to $p_\theta(y_t | y_{<t}, C)$ at each step t , finds a sequence \tilde{Y} whose probability under p_θ is 0 for any context distribution.*



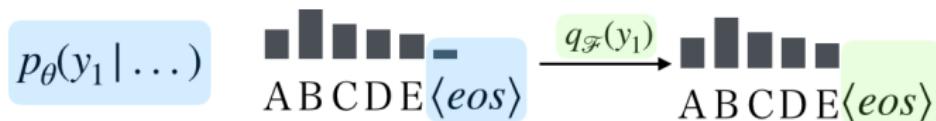
- Is the **model + decoding algorithm** consistent? **No!**

Theorem 3.4 (Inconsistency) of an incomplete decoding algorithm). *There exists a consistent recurrent language model p_θ from which an incomplete decoding algorithm \mathcal{F} , that considers only up to $(|V| - 1)$ -most likely tokens according to $p_\theta(y_t | y_{<t}, C)$ at each step t , finds a sequence \tilde{Y} whose probability under p_θ is 0 for any context distribution.*



- Decoded sequences can be infinite-length!

- ▶ Is the **model + decoding algorithm** consistent? **No!**
- ▶ **Key Idea:** Model might always rank $\langle eos \rangle$ last



RNN always ranks $\langle eos \rangle$ last $\langle eos \rangle$ discarded



- ▶ In summary, we see **infinite sequences** due to a **mismatch** between the **model** p_θ and **decoded** $q_{\mathcal{F}}$ distributions.

- ▶ In summary, we see **infinite sequences** due to a **mismatch** between the **model** p_θ and **decoded** $q_{\mathcal{F}}$ distributions.
- ▶ How do we fix it?

Model

Learning

Inference

Model

Learning

Inference

- ▶ Ensure $\langle eos \rangle$ is eventually sampled

Consistent top- k sampling

Top- k sampling with:

$$q(v) \propto \begin{cases} p_\theta(v | y_{<t}, C), & \text{if } v \in V', \\ 0, & \text{otherwise,} \end{cases}$$

where $V' = \underbrace{\{\langle \text{eos} \rangle\}}_{v'} \cup \arg \max_{v'} p_\theta(v' | y_{<t}, C)$.

Idea: $\langle \text{eos} \rangle$ is always eligible for selection

Consistent top- k sampling

Top- k sampling with:

$$q(v) \propto \begin{cases} p_\theta(v | y_{<t}, C), & \text{if } v \in V', \\ 0, & \text{otherwise,} \end{cases}$$

where $V' = \underbrace{\{\langle \text{eos} \rangle\}}_{v'} \cup \arg \max_{v'} p_\theta(v' | y_{<t}, C)$.

Consistent nucleus sampling

Nucleus sampling with:

$$q(v) \propto \begin{cases} p_\theta(v | y_{<t}, C), & \text{if } v \in V_\mu \cup \underbrace{\{\langle \text{eos} \rangle\}}, \\ 0, & \text{otherwise.} \end{cases}$$

Idea: $\langle \text{eos} \rangle$ is always eligible for selection

Model

Learning

Inference

- ▶ Ensure $\langle eos \rangle$ is eventually top-ranked

Self-terminating Recurrent Language Model

Computes the following conditional probability at each time step:

$$p_{\theta}(v | y_{<t}, C) = \begin{cases} 1 - \alpha(h_t), & \text{if } v = \langle \text{eos} \rangle, \\ \frac{\alpha(h_t) \exp(u_v^\top h_t + c_v)}{\sum_{v' \in V'} \exp(u_{v'}^\top h_t + c_{v'})}, & \text{otherwise,} \end{cases}$$

where

$$\alpha(h_0) = \sigma(u_{\langle \text{eos} \rangle}^\top h_0 + c_{\langle \text{eos} \rangle}),$$

$$\alpha(h_t) = \sigma(u_{\langle \text{eos} \rangle}^\top h_t + c_{\langle \text{eos} \rangle}) [1 - p_{\theta}(\langle \text{eos} \rangle | y_{<t-1}, C)],$$

with $\sigma : \mathbb{R} \rightarrow [0, 1 - \epsilon]$ and $\epsilon \in (0, 1)$.

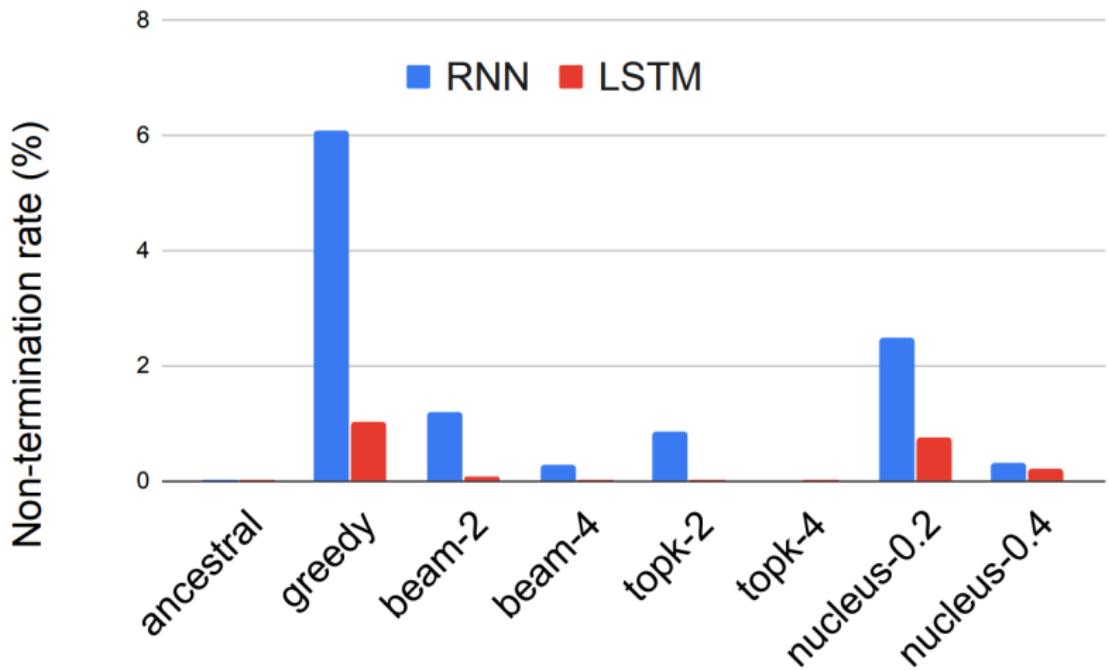
Self-terminating Recurrent Language Model

$$p_t^{\langle \text{eos} \rangle} = 1 - \underbrace{\prod_{t'=0}^t \sigma(a_{t'}^{\langle \text{eos} \rangle})}_{\rightarrow 0}$$

$p_\theta(\langle \text{eos} \rangle | x_{<t}, C)$ increases monotonically.

Does inconsistency occur in practice?

- ▶ Wikitext2; non-termination means no $\langle \text{eos} \rangle$ within 1500 steps



Effect of consistent sampling

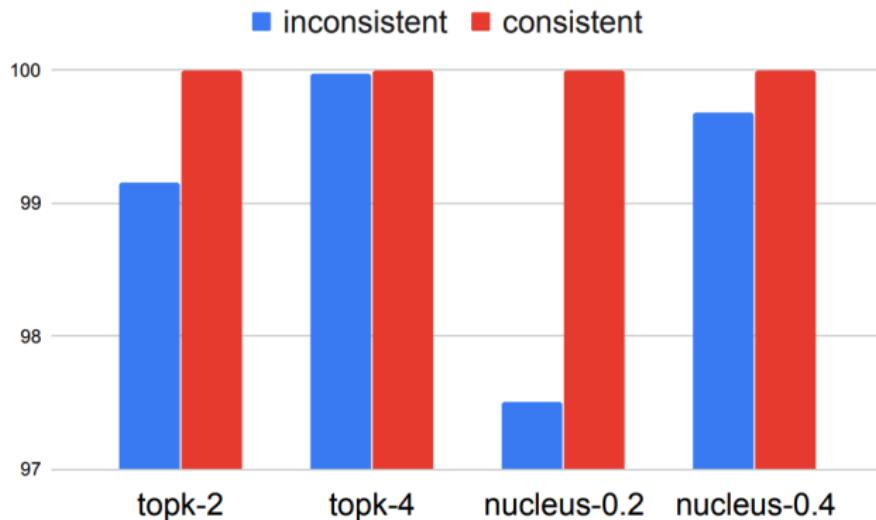


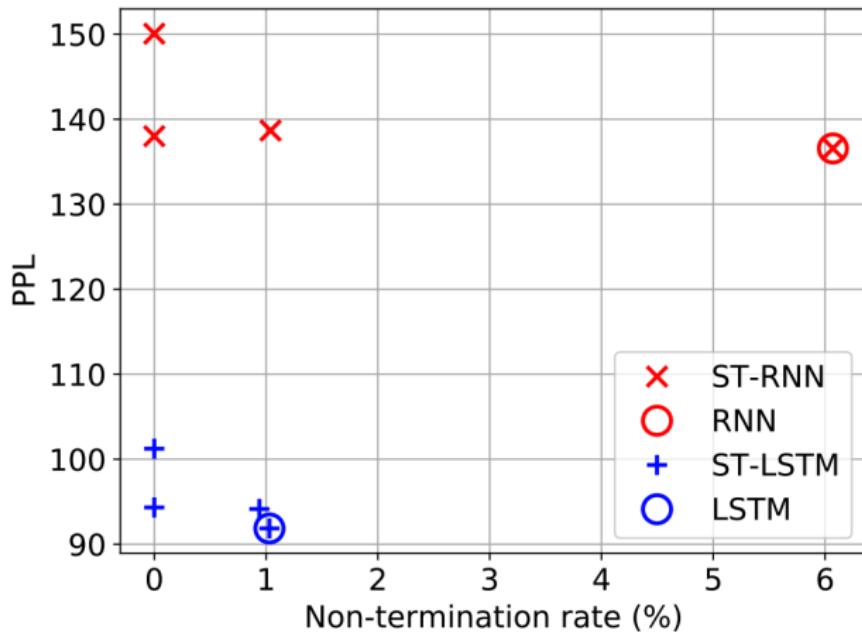
Figure: Termination rate (%) with tanh-RNN

Effect of consistent sampling

Table: Nucleus and consistent nucleus ($\mu = 0.4$) sampling with the LSTM-RNN.

Prefix	<i>He had a guest @-@ starring role on the television</i>
nucleus	film the website , with whom he wrote to the title of The Englishwoman 's Domestic Magazine . ⟨eos⟩
c-nucleus	film the website , but he did not be a new sequel . ⟨eos⟩
Prefix	<i>Somewhere between 29 and 67 species are recognised in the</i>
nucleus	⟨unk⟩ , ⟨unk⟩ → ∞
c-nucleus	⟨unk⟩ , with the exception of an average of 6 @. @ 4 in (1 @. @ 6 mm) . ⟨eos⟩
Prefix	<i>The Civil War saw more ironclads built by both sides</i>
nucleus	and towns , including ⟨unk⟩ , the British Empire , ⟨unk⟩ → ∞
c-nucleus	and towns , including ⟨unk⟩ , the British Empire , ⟨unk⟩ , ⟨unk⟩ , ⟨unk⟩ , ⟨unk⟩ , ⟨unk⟩ , ⟨unk⟩ , ⟨eos⟩

Effect of Self-terminating Recurrent LM



Effect of Self-terminating Recurrent LM

Table: LSTM-RNN and a self-terminating LSTM-RNN ($\epsilon = 10^{-3}$)

Prefix	<i>With 2 : 45 to go in the game ,</i>
Baseline	the team was able to gain a first down . <eos>
STRLM	the Wolfpack was unable to gain a first down . <eos>
Prefix	<i>As of 2012 , she is a horse riding teacher</i>
Baseline	, and a <unk> , → ∞
STRLM	, and a member of the <unk> <unk> . <eos>
Prefix	<i>Nintendo Power said they enjoyed Block Ball and its number</i>
Baseline	of songs , including the " <unk> " , " <unk> " , " <unk> " , ... , " Ode to a Nightingale " , " Ode to a Nightingale " , " Ode to a → ∞
STRLM	of songs , including the " <unk> " , " <unk> " , ... { " <unk> " , } ⁴⁵ ... " <unk> " , " <eos>

Recap

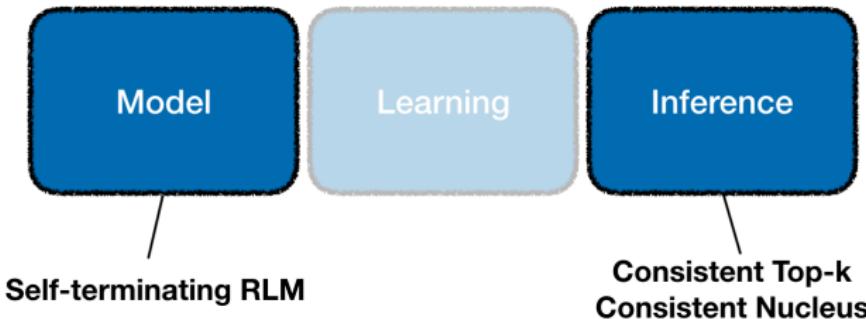
- ▶ Text degeneration is characterized by
non-terminating sequences, repetition, and logical contradictions

Recap

- ▶ Text degeneration is characterized by **non-terminating sequences, repetition, and logical contradictions**
- ▶ Common models + decoding algorithms can yield **infinite sequences** due to ranking $\langle \text{eos} \rangle$ low and excluding it

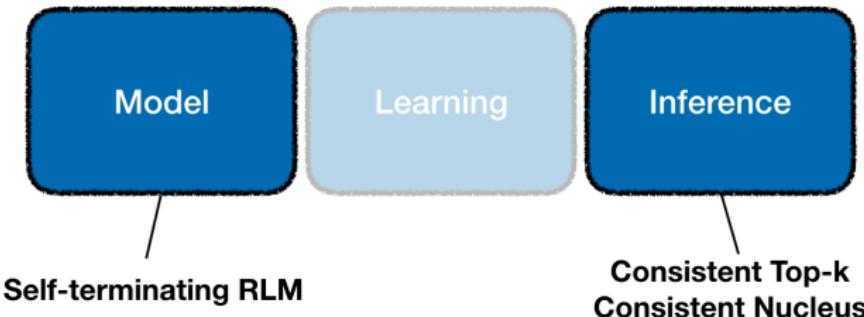
Recap

- ▶ Text degeneration is characterized by **non-terminating sequences, repetition, and logical contradictions**
- ▶ Common models + decoding algorithms can yield **infinite sequences** due to ranking $\langle \text{eos} \rangle$ low and excluding it
- ▶ New **model** and **inference** methods can address the issue



Recap

- ▶ Text degeneration is characterized by **non-terminating sequences, repetition, and logical contradictions**
- ▶ Common models + decoding algorithms can yield **infinite sequences** due to ranking $\langle \text{eos} \rangle$ low and excluding it
- ▶ New **model** and **inference** methods can address the issue



Next: **Repetition & Learning**

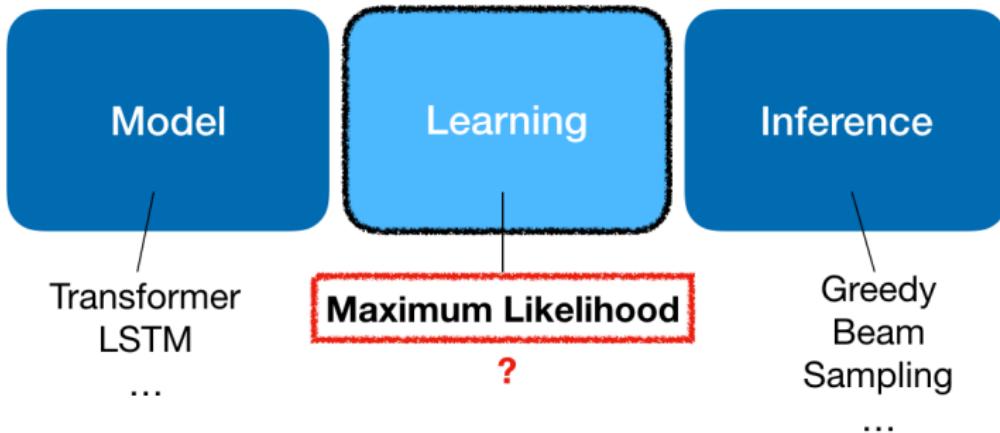
Neural Text Generation with Unlikelihood Training

Welleck, Kulikov, Roller, Dinan, Cho, Weston ICLR 2020



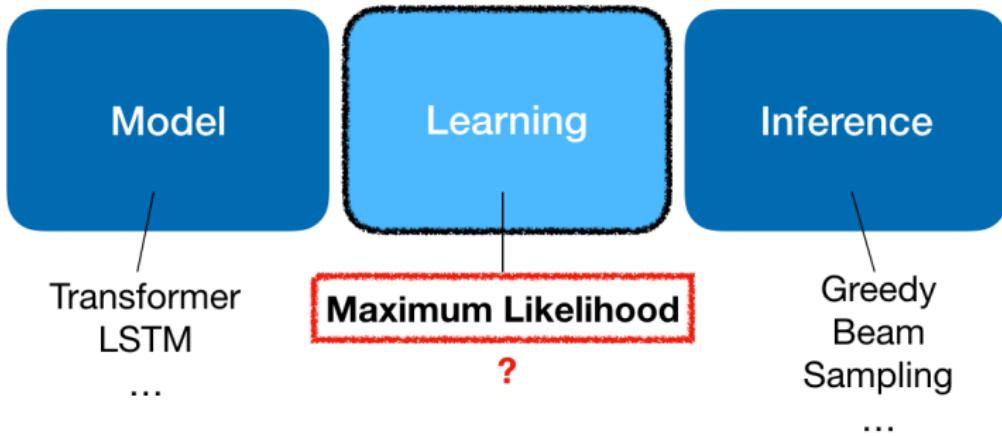
Neural Text Generation with Unlikelihood Training

[Welleck, Kulikov et al ICLR 2020]



Neural Text Generation with Unlikelihood Training

[Welleck, Kulikov et al ICLR 2020]



- ▶ Is there a learning objective that doesn't lead to degeneration?
 - ▶ 1st step: for **specified** degenerate behavior (e.g. repetition).

Maximum Likelihood Estimation (MLE)

Given training example (X, Y) ,

$$\mathcal{L}_{\text{MLE}} = - \sum_{t=1}^{|Y|} \log p_{\theta}(y_t | y_{<t})$$

Maximum Likelihood Estimation (MLE)

Given training example (X, Y) ,

$$\mathcal{L}_{\text{MLE}} = - \sum_{t=1}^{|Y|} \log p_{\theta}(y_t | y_{<t})$$

$$\left. \begin{array}{llll} \text{hello} & \text{how} & \text{are} & \frac{}{y_4} \\ y_1 & y_2 & y_3 & y_4 \end{array} \right\} p_{\theta}(y_4 | y_{<4}, X)$$

a
aardvark
are
...
hello
how
...
yes
you ↑
...
zebra

Maximum Likelihood Estimation (MLE)

Prefix
 \mathcal{L}_{MLE}

... starboard engines and was going to crash . “ We 're going in , ”
he said . “ We 're going to crash . We 're going to
to crash . We 're going to crash . We 're going to crash . We 're going to

Repetition

Model: **40%** 4-gram repetition

Human: **~1%** 4-gram repetition

Large-Scale Transformer LM [from Welleck et al ICLR 2020]

Unlikelihood (Token-level)

$$\mathcal{L}_{ULE} = \mathcal{L}_{MLE} + \alpha \mathcal{L}_{UL}$$

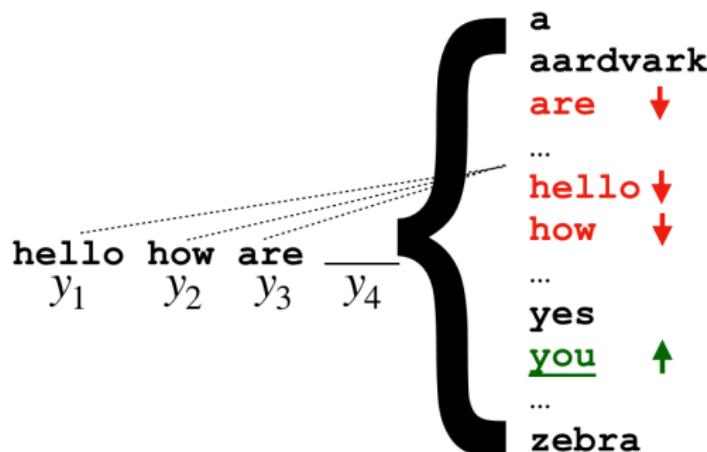
$$\mathcal{L}_{UL}^t = - \sum_{y_{\text{neg}} \in \mathcal{C}_t} \log (1 - p_{\theta}(y_{\text{neg}} | y_{<t}, X))$$

Unlikelihood (Token-level)

$$\mathcal{L}_{ULE} = \mathcal{L}_{MLE} + \alpha \mathcal{L}_{UL}$$

$$\mathcal{L}_{UL}^t = - \sum_{y_{\text{neg}} \in \mathcal{C}_t} \log (1 - p_{\theta}(y_{\text{neg}} | y_{<t}, X))$$

$$\mathcal{C}_t = \{y_1, \dots, y_{t-1}\} \setminus \{y_t\}$$



Unlikelihood (Sequence-level)

$$\hat{Y} \sim \mathcal{F}(p_\theta, X)$$

decoded sequence

$$\mathcal{L}_{ULE} = \mathcal{L}_{MLE}(X, Y) + \alpha \mathcal{L}_{UL}(X, \hat{Y}, \mathcal{C})$$

Unlikelihood (Sequence-level)

$$\hat{Y} \sim \mathcal{F}(p_\theta, X)$$

decoded sequence

$$\mathcal{L}_{ULE} = \mathcal{L}_{MLE}(X, Y) + \alpha \mathcal{L}_{UL}(X, \hat{Y}, \mathcal{C})$$

Candidates \mathcal{C} :

- ▶ Random $\mathcal{C}_t = \{\hat{y}_t\}$ if `rand()` else \emptyset

he said . we're going to crash . we're going to crash

\hat{Y}

Unlikelihood (Sequence-level)

$$\hat{Y} \sim \mathcal{F}(p_\theta, X)$$

decoded sequence

$$\mathcal{L}_{ULE} = \mathcal{L}_{MLE}(X, Y) + \alpha \mathcal{L}_{UL}(X, \hat{Y}, \mathcal{C})$$

Candidates \mathcal{C} :

- ▶ Repeats $\mathcal{C}_t = \{\hat{y}_t\}$ if $\hat{y}_t \in \text{repeat}$, else \emptyset

he said . we're going to crash . we're going to crash

\hat{Y}

Unlikelihood (Sequence-level)

$$\hat{Y} \sim \mathcal{F}(p_\theta, X)$$

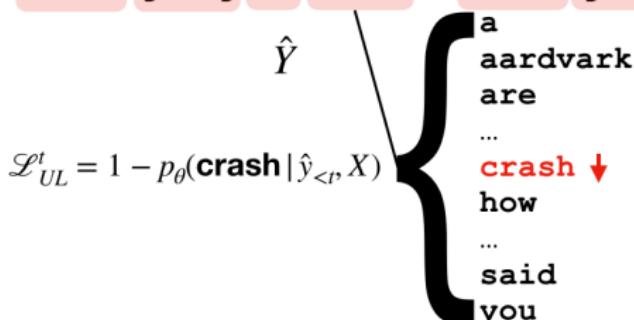
decoded sequence

$$\mathcal{L}_{ULE} = \mathcal{L}_{MLE}(X, Y) + \alpha \mathcal{L}_{UL}(X, \hat{Y}, \mathcal{C})$$

Candidates \mathcal{C} :

- ▶ Repeats $\mathcal{C}_t = \{\hat{y}_t\}$ if $\hat{y}_t \in \text{repeat}$, else \emptyset

he said . we're going to crash . we're going to crash



Model	search	seq-rep-4	uniq-seq	ppl
MLE	greedy	.442	10.8k	25.64
	beam	.523	9.5k	
UL - Random	greedy	.116	12.5k	25.518
	beam	.146	14.2k	
UL - Repeat	greedy	.058	15.4k	26.72
	beam	.013	19.1k	
Human	-	.006	19.8k	-

Winner		Loser	Crowdworkers	Experts
			Win rate	Win rate
Unlikelihood	<i>beats</i>	MLE	*82%	
Unlikelihood	<i>beats</i>	MLE + Nucleus Sampling (0.9)	59%	*83%
Unlikelihood		MLE + Beam Blocking (4 gram)	60%	*74%

Why does it work?

$$\nabla \mathcal{L}_a = x^* - m \odot p, \quad m_i = \begin{cases} (1 - \alpha \frac{p_{\text{neg}}}{1-p_{\text{neg}}}) & \text{if } i \neq i_{\text{neg}} \\ (1 + \alpha) & \text{if } i = i_{\text{neg}}, \end{cases}$$

1. **Negative candidate:** negative gradient
2. **Ground-truth token:** positive gradient
3. **Others:** positive when negative candidate probability is high

Why does it work?

- ▶ Noise contrastive estimation [Gutmann & Hyvonen 2010]

$$\underbrace{\mathbb{E}_{x_+ \sim p_+} [\log h_\theta(x_+)]}_{\textit{positives}} + \underbrace{\mathbb{E}_{x_- \sim p_-} [\log(1 - h_\theta(x_-))]}_{\textit{negatives}}$$

Where

$$h_\theta(x_-) = \sigma(\log p_\theta(x) - \log p_-(x))$$

Why does it work?

- ▶ Noise contrastive estimation [Gutmann & Hyvonen 2010]

$$\underbrace{\mathbb{E}_{x_+ \sim p_+} [\log h_\theta(x_+)]}_{\text{positives}} + \underbrace{\mathbb{E}_{x_- \sim p_-} [\log(1 - h_\theta(x_-))]}_{\text{negatives}}$$

Where

$$h_\theta(x_-) = \sigma(\log p_\theta(x) - \log p_-(x))$$

- ▶ Unlikelihood with random candidates:

$$\underbrace{\mathbb{E}_{\substack{x \sim p_+ \\ t \sim \mathcal{U}(T)}} [\log p_\theta(x_t | x_{<t})]}_{\text{positives}} + \underbrace{\mathbb{E}_{\substack{x \sim p_- \\ t \sim \mathcal{U}(T)}} [\log(1 - p_\theta(x_t | x_{<t}))]}_{\text{negatives}}.$$

Follow-up: Unlikelihood for Dialogue

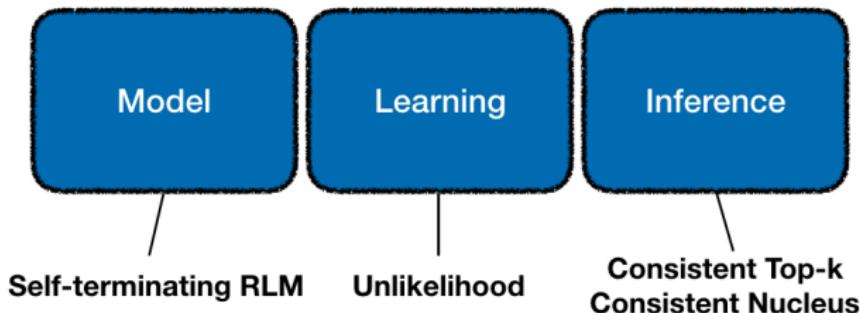
Don't Say That!
Making Inconsistent Dialogue Unlikely with Unlikelihood Training

**Margaret Li, Stephen Roller, Ilia Kulikov, Sean Welleck
Y-Lan Boureau, Kyunghyun Cho, Jason Weston**
Facebook AI Research

- ▶ Repetition and logical contradictions

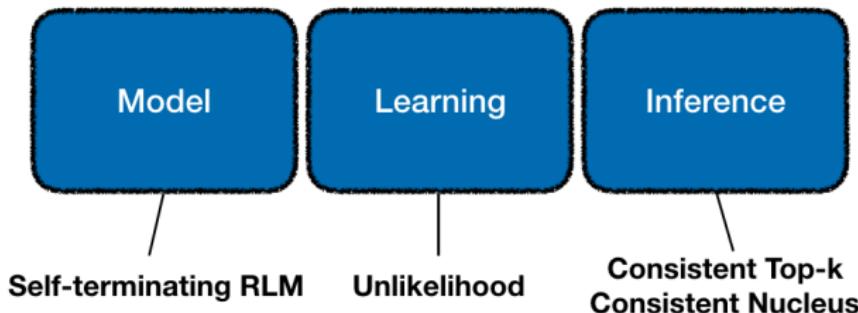
Recap

- ▶ Text degeneration is characterized by **non-terminating sequences, repetition, and logical contradictions**
- ▶ New **model** and **inference** methods can address the issues



Recap

- ▶ Text degeneration is characterized by **non-terminating sequences, repetition, and logical contradictions**
- ▶ New **model** and **inference** methods can address the issues



But more investigation and methods are needed!

**Thank you
very much very much very much very much very → ∞**